

# Machine Learning (Theory)

15/1/23

cls-1

→ Defn lab 275 wr,

Machine Learning:

- \* A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience — Tom Michalek (1998)

- \* Arthur Samuel (1959):

## Machine Learning

Supervised Learning

when input, output  
given & feature  
leveled.

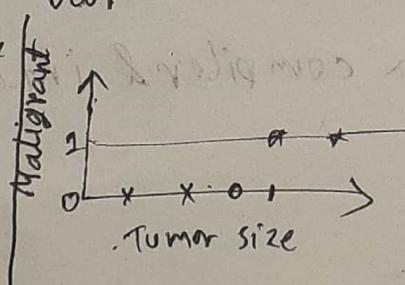
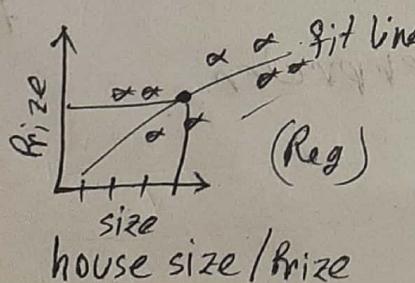
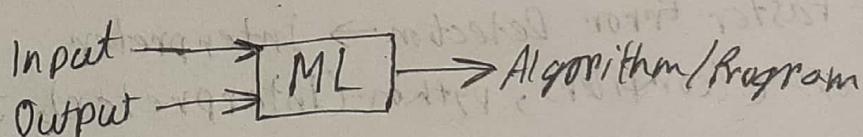
Unsupervised Learning

No output, only  
input & feature  
not leveled.

Reinforcement Learning

when machine learns  
through trial & errors

- \* Supervised learning:



(cls) → Yes or no  
page  
size

finding

Supervised learning sets:-

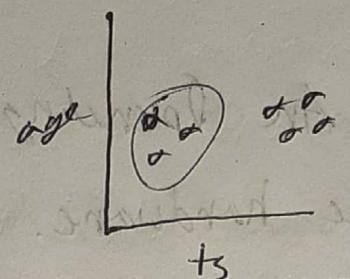
- Regression (Continuous) - value fitting
- Classification (Discrete) - given class / identify wr

# Supervised Learning :-

Input ( $x$ )	Output ( $y$ )
1	5
3	10
5	15
10	25

Model/H/F

# Unsupervised :-



# Reinforcement :-

Action  $\rightarrow$  against reward point 2000.Regression :- output space continuousClassification :-      in      discrete\* Regression:Linear Regression  $\rightarrow$  straight line drawn for\* Training Data:- (House size & Price)

$x$	size ( $ft^2$ )	Price (k)
2000	400	
1500	350	
1200	280	
500	150	
:	:	
$n = 50$		

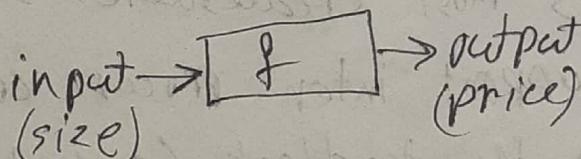
 $x$  = input $y$  = output $n$  = number of data $(x^i, y^i)$  =  $i^{th}$  training example $(x^3, y^3) = (1200, 280)$  $x$  = input,  $\hat{y}$  = prediction

training set

Learning Algorithm

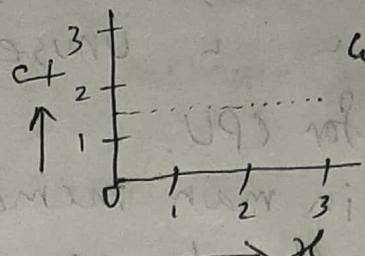
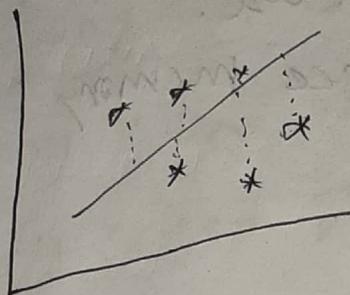
function (output)

hypothesis / model / function



\* Linear Regression:  $f = w\mathbf{x} + b$  [ $w$  = slope,  $b$  = constant]

$w, b$ : parameters, so they are changeable



$$w = 0, b = 1.5$$

\* Cost function:  $J_{w,b} = w\mathbf{x} + b$

Squared error cost function:-

$$J_{w,b} = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

$$J_{w,b} = \frac{1}{2n} \sum_{i=1}^n (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Objective:- to minimize cost function ( $J$ ).

\* Task:  $f_{w,b} = w\mathbf{x}$ , where  $b=0$ ,  $\frac{\partial J}{\partial w} = \frac{1}{2n} \sum_{i=1}^n (f_{w,b}(x^{(i)}) - y^{(i)})^2$

$w=1$ :  $J_1 = \frac{1}{2 \times 3} (1-1)^2 + (2-2)^2 + (3-3)^2 = 0$

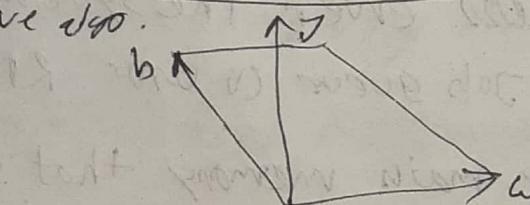
$w=0.5$ :  $J_{0.5} = \frac{1}{2 \times 3} (0.5-1)^2 + (1-2)^2 + (1.5-3)^2 = 0.58$

$w=1.5$ :  $J_{1.5} = \frac{1}{2 \times 3} (1.5-1)^2 + (3-2)^2 + (4.5-3)^2 = 0.58$

$w=0$ :  $J_0 = \frac{1}{2 \times 3} (0-1)^2 + (0-2)^2 + (0-3)^2 = 2.33$

\* Gradient Descent Algorithm: used for cost/loss function optimization.

Iterative algo.

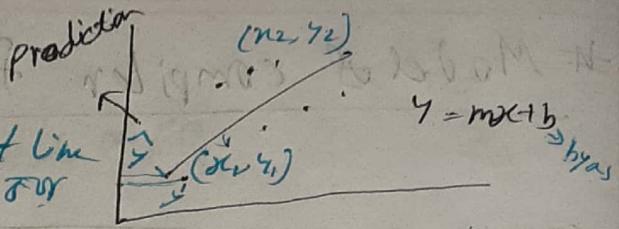


$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial}{\partial w} J(w, b)$$

learning rate

$$b_{\text{new}} = b_{\text{old}} - \alpha \frac{\partial}{\partial b} J(w, b)$$

## # Linear Regression :-

→ Squared Error: calculated best fit line  
(SSE)

$$\text{minimization (goal)} \quad SE = \sum_{i=1}^n (y_i - (mx_i + b))^2$$

$$= y_1^2 + 2y_1(mx_1 + b) + (mx_1 + b)^2 + y_2^2 + 2y_2(mx_2 + b) + (mx_2 + b)^2 + \dots + y_n^2 + 2y_n(mx_n + b) + (mx_n + b)^2$$

$(m, b \text{ are constants})$   
Simplify by  $y^2$

$$SE = (y_1^2 + y_2^2 + \dots + y_n^2) - 2m(x_1y_1 + x_2y_2 + \dots + x_ny_n) - 2b(y_1 + y_2 + \dots + y_n) + m^2(x_1^2 + x_2^2 + \dots + x_n^2) + 2mb(x_1 + x_2 + \dots + x_n) - nb^2$$

Here,  $\frac{y_1^2 + y_2^2 + \dots + y_n^2}{n} = \bar{y}^2$  (average of  $y^2$ )

$$= y_1^2 + y_2^2 + \dots + y_n^2 = n\bar{y}^2 \quad \text{--- } ①$$

again,

$$x_1y_1 + x_2y_2 + \dots + x_ny_n = n\bar{x}\bar{y} \quad \text{--- } ②$$

$$SE = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2b\bar{n}\bar{y} + m^2\bar{n}^2 + 2mb\bar{n}\bar{x} + nb^2$$

At minimum SE:  $\frac{\delta SE}{\delta m} = 0$

$$\Rightarrow 0 - 2n\bar{x}\bar{y} + 2n\bar{x}^2m + 2b\bar{n}\bar{x} = 0$$

[In for divide]  $\Rightarrow m\bar{n}^2 + b\bar{n} = \bar{y} \quad \text{--- } ③$

$$\frac{\text{SSE}}{n} = 0$$

$$\Rightarrow -2n\bar{y} + 2nx^T m + 2bn\bar{x} = 0$$

$$\Rightarrow m\bar{x} + b = \bar{y} \quad \textcircled{2}$$

$$\Rightarrow m = \frac{\bar{y} - \bar{b}\bar{x}}{(\bar{x})^2 - \bar{x}^2} = \frac{\bar{y} - \bar{b}\bar{x}}{n^2 - (\bar{x})^2}$$

$$b = \bar{y} - m\bar{x}$$

\* Gradient decent Algorithm: (Feature-Size & Data-Size)

$$\begin{aligned} f(n) &= y = mx + b = \underset{\substack{\uparrow \\ \text{independent}}}{w^T x} + b \quad h_{\theta}(x) = \theta_0 + \theta_1 x_1 \\ &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\ h(\theta) &= \sum_{j=0}^n \theta_j x_j = \theta^T x \quad \left[ \begin{matrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_n \end{matrix} \right] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \end{aligned}$$

\* Cost function  $J(\theta) = \frac{1}{2n} \sum_{i=1}^n \frac{(h_{\theta}(x^{(i)}) - y^{(i)})^2}{\text{prediction} - \text{true}}$

Objective:  $\theta^*$  = argument minimum  $J(\theta)$

ADA

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} - \frac{\partial J(\theta)}{\partial \theta_j}$$

learning rate / rate of change / slope

## # Batch Gradient Decent

Full data at a time

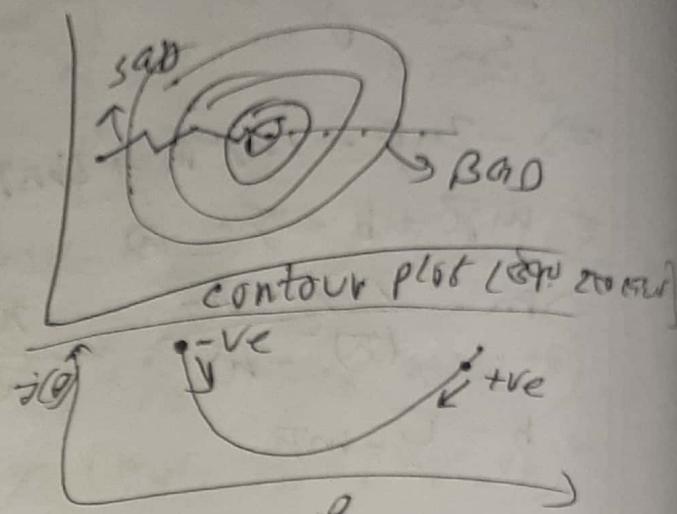
~~slow~~

## # Stochastic Gradient Decent

At a time one update

## # Minibatch gradient descent

For single Points



$$\begin{aligned}
 n=1 \quad \frac{\partial}{\partial \theta_j} j(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\
 &= \frac{1}{2} \cdot 2 (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\
 &= (h_\theta(x) - y) n_j \\
 \therefore \theta_j^{old} - \alpha (h_\theta(x) - y) n_j &\quad \left. \begin{array}{l} (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_j x_j + \dots + \theta_n x_n) \\ + \theta_j n_j + \dots + \theta_n x_n \end{array} \right)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \theta_j x_i - y &= \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_j x_j + \dots + \theta_n x_n - y) \\
 \frac{\partial \theta_j n_j}{\partial \theta_j} &= n_j
 \end{aligned}$$

Repeat until converge  $\rightarrow$  GDA

for  $i = 1$  to  $n$ :

$$\theta_j = \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) / n, \text{ (i)} \quad \text{which nb feature}$$

$\downarrow$  G.D.A.

$x^{(i)}, y^{(i)}$  = input data

# ML (T-4)

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = \frac{\theta^T}{s} \frac{X}{\sqrt{s}}$$

[ $\vec{\theta}$ -d vector =  
 $\vec{x}$ -d feature]

Loss function:  $J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$

## Probabilistic Interpretation

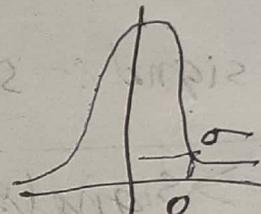
$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

error/noise

$$\epsilon^{(i)} \sim N(0, \sigma^2)$$

mean variance

$$\begin{cases} x^{(i)} \in \mathbb{R}^{d+1} \\ y^{(i)} \in \mathbb{R} \end{cases}$$



$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

Gaussian distribution

$$p(y^{(i)} | x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

given  $x^{(i)}$  fixed

## Maximum Likelihood Estimation

Likelihood function:

$$\begin{aligned} L(\theta) &= p(\vec{y} | \vec{x}; \theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

log

Log Likelihood:-

$$\begin{aligned}
 l(\theta) &= \log L(\theta) \\
 &= -\log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
 &= -\sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
 &= -n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2
 \end{aligned}$$

→ likelihood  $\propto$  maximum joint probability

Hence, maximizing  $l(\theta)$  gives the same answer as minimizing  $\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = J(\theta)$

## # Classification

$y \in \{0, 1\}$  binary classification

$$h_\theta(x) \in [0, 1]$$

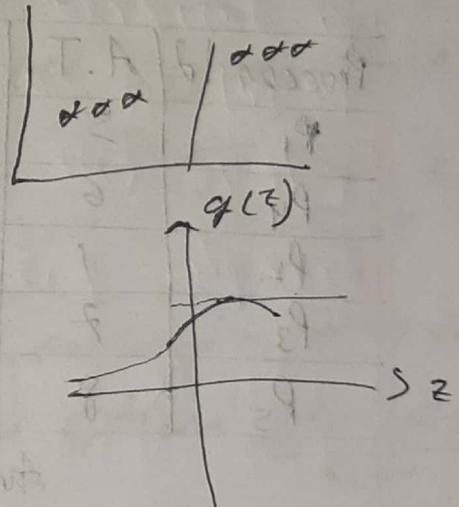
$$h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1 + e^{-z}} \text{ (sigmoid function)}$$

$$P(y=1/x; \theta) = h_\theta(x)$$

$$P(y=0/x; \theta) = 1 - h_\theta(x)$$

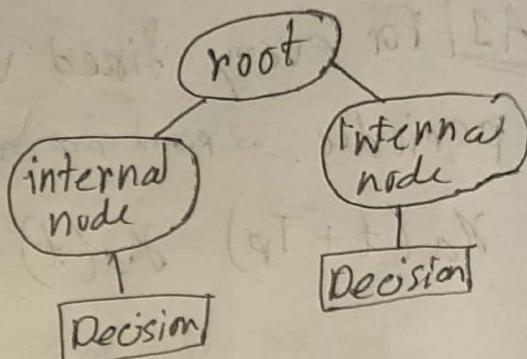
$$P(y/x; \theta) = h_\theta(x)^y (1 - h_\theta(x))^{1-y}$$



ML (2-3)# Decision Tree:

Q. Tennis (car 20) for?

Day	Outlook	Temp	Humidity	Tennis
D <sub>1</sub>	Sunny	Hot	High	No
D <sub>2</sub>	Sunny	Mild	Normal	Yes
D <sub>3</sub>	Rain	Cold	Normal	No



\* Information Gain (IG) :-

$$IG = \text{Entropy}(S) - \sum_{i=1}^c \frac{|S_i|}{|S|} \text{Entropy}(S_i)$$

\* Entropy :- It is a measure from information theory, ~~characteristics~~ characterize the (im)purity/homogeneity of an arbitrary collection.

$$E = -\sum_{i=1}^c p_i \log_2(p_i)$$

\* Factors affecting sunburn:-

Name	Hair	Height	Weight	Lotion	Decision
Sara	Blonde	avg	light	No	✓
Dana	Blonde	tall	avg	Yes	✗
Alex	Brown	short	avg	No	✗
Anny	Blonde	short	avg	No	✗
Emily	Red	avg	heavy	No	✗
Pete	Brown	tall	heavy	No	✗
Jhon	Brown	avg	heavy	No	✗
Kati	Blonde	short	light	Yes	✗

$$0.5 - 0.2175$$

$$\frac{S_i}{S} \times \left( - \sum p_i \log_2(p_i) \right)$$

$$\rightarrow \text{avg entropy for hair color} = \frac{4}{8} \times \left\{ - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right\}$$

Blond

$$+ \frac{3}{8} \times \left\{ - \frac{3}{3} \log_2 \left( \frac{3}{3} \right) - \frac{0}{3} \right\} + \frac{1}{8} \times \left\{ - \frac{1}{1} \log_2 \left( \frac{1}{1} \right) - \frac{0}{1} \right\}$$

Brown

$$= \underline{0.5}$$

Non

Red

Non

$$\rightarrow \text{avg entropy for height} = 2 \times \left[ \frac{3}{8} \times \left\{ - \frac{2}{3} \log_2 \left( \frac{2}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right\} + \right.$$

$$\left. \frac{2}{8} \times \left\{ - \frac{2}{2} \log_2 \left( \frac{2}{2} \right) - \frac{0}{2} \right\} \right]$$

avg, short

tall

$$= \underline{0.69}$$

$$\rightarrow \text{avg entropy for weight} = \frac{2}{8} \left\{ - \frac{1}{2} \log_2 \left( \frac{1}{2} \right) - \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right\} +$$

$$2 \times \left[ \frac{3}{8} \left\{ - \frac{2}{3} \log_2 \left( \frac{2}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right\} \right]$$

avg, heavy

light

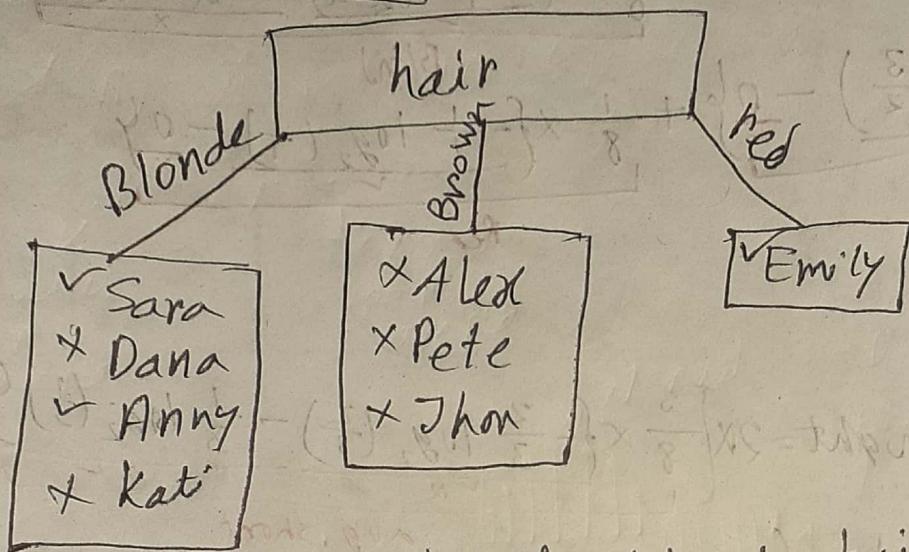
$$= \underline{0.936}$$

$$\rightarrow \text{avg entropy for lotim} = \frac{2}{8} \left\{ - \frac{2}{2} \log_2 \left( \frac{2}{2} \right) \right\} + \frac{6}{8} \left\{ - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right\}$$

$$= \underline{0.75}$$

$$\rightarrow \text{Decision } (\epsilon_S) = \frac{3}{8} \log_2 \left( \frac{3}{8} \right) - \frac{5}{8} \log_2 \left( \frac{5}{8} \right)$$

Information Gain:  $E(S) - 0.5$

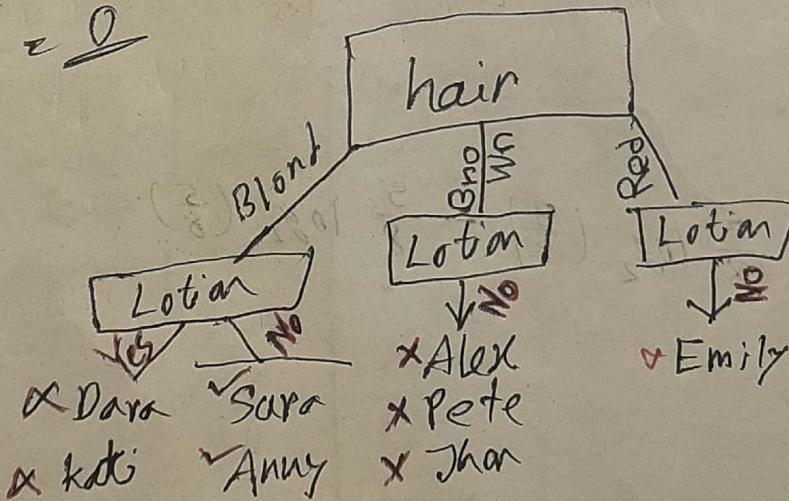


Another test needed for blonde hair color :-

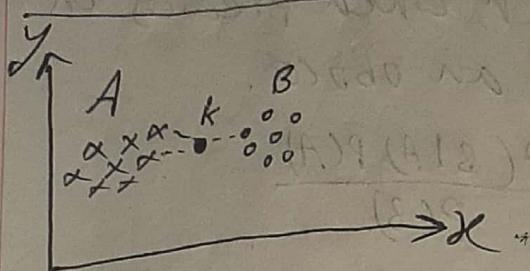
$$\text{Height} = \frac{2}{4} \left\{ -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right\} + 0 + 0$$

$$\text{Weight} = 2 \times \left[ \frac{2}{4} \left\{ -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right\} + 0 + 0 \right]$$

$$\text{Lotion} = 2 \times \left[ \frac{2}{4} \left\{ -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \log_2 (-) \right\} + 0 \right]$$



## # K-Nearest Neighbors Algorithm (KNN):



→ k select from 20 1st

→ If data point 20 distance calculate from 20

→ distance sort

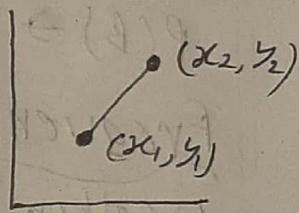
→ Nearest data point can set = 20 from 20  
(nb of values close to k)

- Parametric Algorithm:- learning to accept w thru zero loss
- Non-parametric Algorithm:- Ex:- Decision Tree, KNN
- k tips:- (i) odd number, (ii) odd feature, (iii) medium value.

## # Distance Measure:-

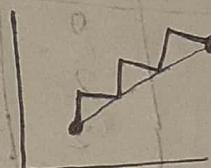
→ Euclidean distance:- for points 20 20

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

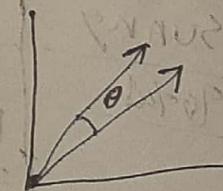


→ Manhattan distance:-

$$|x_1 - x_2| + |y_1 - y_2|$$



→ Cosine similarity =  $\frac{|A||B| \cos \theta}{|A||B|}$  [self study]



### Advantage of KNN

- Easy to implement.
- No need to build a model, tune several parameters or make additional assumption.
- It can be used for classification as well as regression.

### Disadvantage of KNN

- The algorithm gets significantly slower as the number of examples (training data) and/or predictions (test data) increase.
- Memory requirement is high.

\* Application:- Decision Tree, Recommended system, Data processing, classification.

## # Naive Bayes Algorithm:-

It is a probabilistic classifier, which predicts on the basis of the probability of an object.

Bayes theorem:-  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

where,

$P(A|B)$  → Probability of hypothesis A on the observed event B  
 Posterior probability

$P(B|A)$  → Probability of hypothesis A given the evidence given likelihood that the probability of hypothesis is true.

$P(A)$  → Probability of hypothesis before observing the evidence

$P(B)$  → Probability of evidence.

### Frequency table:

Weather	Y (Yes)	N (No)	Probability
Overcast	5	0	$5/14 = 0.357$
Rainy	2	2	$4/14 = 0.29$
Sunny	9	2	$5/14 = 0.357$
Total	16	4	14
	$10/14$	$4/14$	
	$= 0.71$	$= 0.28$	

Q. if weather = sunny, then play = Yes / No?

$$P(S|Y) = \frac{3}{10} = 0.3$$

$$P(S) = \frac{5}{14} = 0.35$$

$$P(Y) = \frac{10}{14} = 0.71$$

$$P(S|N) = \frac{2}{4} = 0.5$$

$$P(N) = \frac{4}{14} = 0.28$$

Weather	Play
Rainy	Y
Sunny	Y
Overcast	Y
Overcast	Y
Sunny	N
Rainy	Y
Sunny	Y
Overcast	Y
Rainy	N
Sunny	N
Sunny	Y
Rainy	N
Overcast	Y
Overcast	Y

$$P(Yes|Sunny) = \frac{P(Sunny|Yes) * P(Yes)}{P(Sunny)}$$

$$= \frac{0.3 * 0.71}{0.35} = 0.6$$

$$P(No|Sunny) = \frac{P(Sunny|No) * P(No)}{P(Sunny)} = \frac{0.5 * 0.28}{0.35} = 0.408$$

since  $P(Yes|Sunny) > P(No|Sunny)$  therefore decision = Yes.